

line 1 - B

---

**MANE 3332.01**

# LECTURE 18

# Agenda

- Midterm exams are not graded; still contacting students who missed
- Linear Combination Practice Problems (assigned 10/28, due 10/30)
- Linea Combination Quiz (assigned 10/30, due 11/4)
- Complete Chapter Six and Start Chapter 7
- Attendance
- Questions?

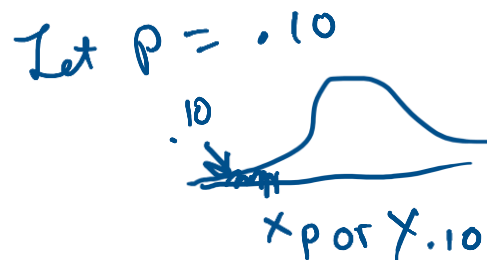
# Handouts

- Lecture 18 Slides
- Lecture 18 Slides - marked

## **CHAPTER 6, CONTINUED**

# Find $x$ such that $P(X \leq x) = p$ Calculating Quantiles

$$P(X \leq x) = p$$



Quartiles

$$p = .25 \text{ or } .5 \text{ or } .75$$

Deciles

$$p = .1 \text{ or } .2 \text{ or } \dots .9$$

Percentiles

$$p = .01, .02, \dots .99$$

reference for calculating quantiles

## 2.3.2 Sample Quantiles

In Example 2.8, we consider an ogive for the plated bracket data. The point  $(1.55, 0.567)$  is on that ogive, so we estimate that 56.7% of the sampled population of brackets weighed at most 1.55 ounces. Weights associated with other percentages can also be estimated by locating the appropriate point on the ogive. In general, if the point  $(x, p)$  is on the ogive, we can use  $x$  as an estimate of the weight with 100% of the population values at or below it. This estimate, called the 100th sample quantile, is denoted  $x_p$ .

If two persons (or computer programs) use different groupings to obtain an ogive, the resulting quantiles will differ. To remedy this deficiency, an algebraic procedure is required.

### THE 100th SAMPLE QUANTILE

Several definitions of sample quantiles are used. We use the one that agrees with the default values output by the UNIVARIATE procedure in SAS\*. Also, the definition used here is consistent with our definition of the sample median.

Suppose a sample of size  $n$  is obtained from some population associated with a continuous variable. For  $0 < p < 1$ , let  $p(n+1)$ ,  $i = \lfloor p(n+1) \rfloor$ , with  $i$  the integer part of  $p(n+1)$  and  $0 \leq d < 1$  the decimal part. If  $1 \leq i < n$  and  $d = 0$ , the 100th sample quantile is  $x_{(i)}$ . If  $1 \leq i < n$  and  $0 < d < 1$ , interpolate linearly between  $x_{(i)}$  and  $x_{(i+1)}$ . In either case, the 100th sample quantile is

$$x_p = x_{(i)} + d[x_{(i+1)} - x_{(i)}] \quad (2.4)$$

when  $1 \leq i < n$ . If  $n = 0$  or  $n$ , the 100th sample quantile does not exist. If 100% is an integer, the corresponding quantile is called a *percentile*.

### EXAMPLE 2.18

Suppose we want to find the 45th percentile of the sample of plated weights in Table 2.1. Since

there are  $n = 75$  observations in the sample and  $p = 0.45$ , we find  $p(n+1) = (0.45)(75+1) = 32.68$ . Letting  $i = 32$  and  $d = 0.68$ , we use Equation (2.4) to obtain  $x_{.45} = x_{(32)} + (0.68)(x_{(33)} - x_{(32)})$ . The 32nd ordered value in Figure 2.1(b) is  $x_{(32)} = 1.50$  and the 33rd ordered value is  $x_{(33)} = 1.51$ . Thus, the 45th percentile for these data is  $x_{.45} = 1.50 + (0.68)(1.51 - 1.50) = 1.5068 \approx 1.507$ . Using this as a point estimate of the population percentile, we can say that approximately 45% of the plated brackets produced on the day the data were collected had weights of 1.507 ounces or less.

The Sample Median is a Percentile. Suppose we want to find the 50th percentile and the data set contains  $n$  values. When  $n$  is even,  $(0.50)(n+1) = (n/2) + (0.50)$ , with  $n/2$  a positive integer. Using Equation (2.4) with  $i = n/2$  and  $d = 0.50$ ,  $x_{.50} = x_{(i)} + (0.50)[x_{(i+1)} - x_{(i)}] = [x_{(i)} + x_{(i+1)}]/2$ . When  $n$  is odd,  $(0.50)(n+1) = (n+1)/2$ , with  $(n+1)/2$  a positive integer. Using Equation (2.4) with  $i = (n+1)/2$  and  $d = 0$ , we find  $x_{.50} = x_{(i)}$ . But, this is precisely how the sample median was defined. Thus,  $x = x_{.50}$ .

### SAMPLE QUANTILES

The percentiles  $x_{.25}$ ,  $x_{.50}$ , and  $x_{.75}$  are known as the *first*, *second*, and *third sample quantiles*, respectively. These quantiles are often denoted  $q_1$ ,  $q_2$ , and  $q_3$ .

### EXAMPLE 2.19

Consider the plated bracket weights in Table 2.1. Using the ordered stem-and-leaf display presented in Figure 2.1(b), we find the following:

- First Quantile:* Since  $(0.25)(75+1) = 19$ ,  $q_1 = x_{.25} = x_{(19)} = 1.46$ .
- Second Quantile (Median):* Since  $(0.50)(75+1) = 38$ ,  $q_2 = x = x_{.50} = x_{(38)} = 1.53$ .

1, 2, 3, 4 median  
 $p = .5$   
 $j(4n) = 2.5$   
 $i = 2, d = .5$   
 $x_p = x_{(i)} + d[x_{(i+1)} - x_{(i)}]$

Quantiles

# Quantile Example

$x_{(i)}$  - ~~with~~ rank observations

8 Observations from binomial distribution with  $n = 10$  and  $p = .5$

	1	2	3	4	5	6	7	8
$x_i$	6	4	5	7	3	5	4	6

$x_{(i)}$  3 4 4 5 5 6 6 7

$$x_2 = 4, x_{(2)} = 4, x_3 = 5, x_{(3)} = 4$$

$$x_p = x_{(i)} + d[x_{(i+1)} - x_{(i)}]$$

$$p = .75$$

$$p_{(n+1)} = .75(8+1) = \underline{6.75}$$

$1 = 6 \text{ \& } d = .75$

## Quantile Example

$$x_{.75} = x_{(6)} + .75[x_{(7)} - x_{(6)}]$$

$$= 6 + .75(6 - 6)$$

# Exploratory Data (Graphical) Analysis

- Exploratory data analysis (EDA) is the use of graphical procedures to analyze data.
- John Tukey was a pioneer in this field and invented several of the procedures
- Tools include stem-and-leaf diagrams, box plots, time series plots and digidot plots



# Stem and Leaf Diagram

- Excellent tool that maintains data integrity
- The stem is the leading digit or digits
- The leaf is the remaining digit
- Make sure to include units
- R Code

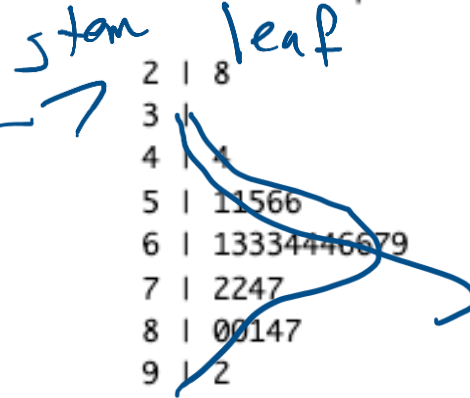
```
stem(midterm$MidtermExam)
```

leaf unit

### Stem and Leaf Example

R output of a Stem and Leaf diagram

Stem = 2  
leaf = 8  
28.0



The decimal point is 1 digit(s) to the right of the |

Stem and Leaf Plot of Midterm Exam Scores

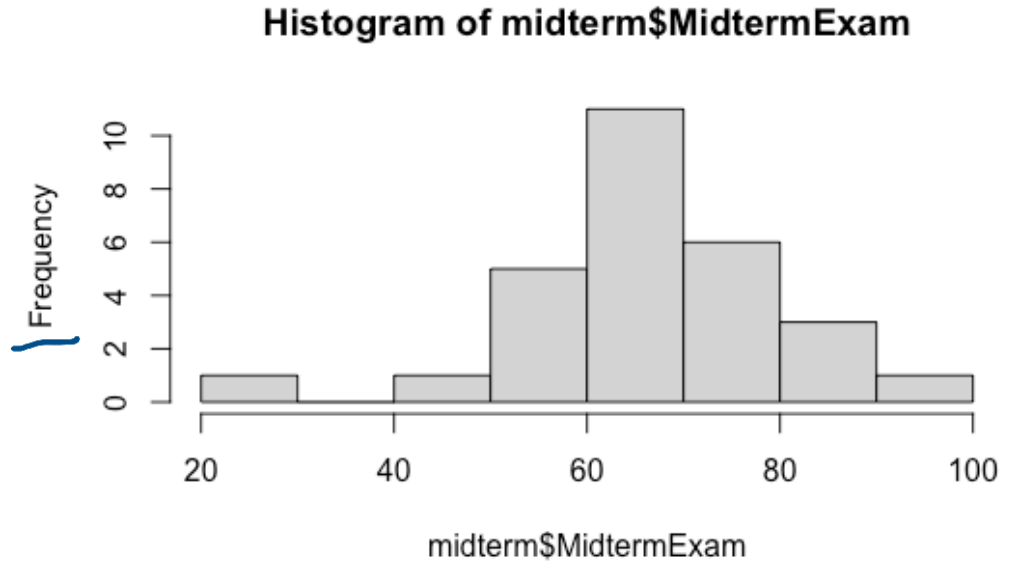
# Histogram

- A histogram is a barchart displaying the frequency distribution information
- There are three types of histograms: frequency, relative frequency and cumulative relative frequency
- R code

```
hist (midterm$MidtermExam)
```

Histogram Example  
R output of histogram

Relative  
Frequency



Histogram of Midterm Exam Scores

# Box & Whiskers

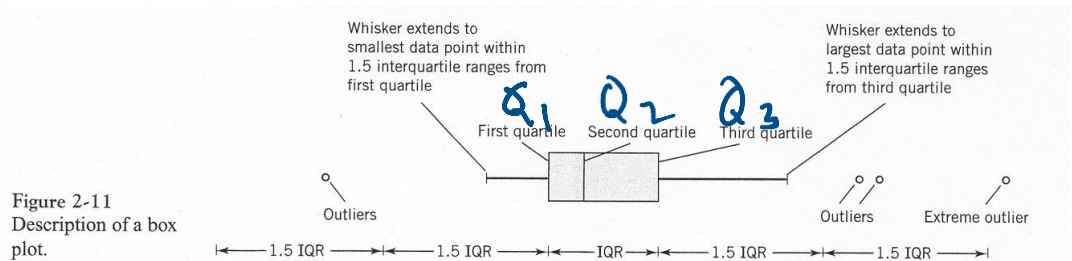
## Boxplot

Graphical display that simultaneously describes several important features of a data set such as center, spread, departure from symmetry and outliers

Requires the calculation of quantiles (quartiles)

### Box Plot 1

IQR - inter quartile range  
 $= Q_3 - Q_1$



potential outliers

Box plot with explanation

# Box Plot 2

*Highest Quality level*  
*Most variability*

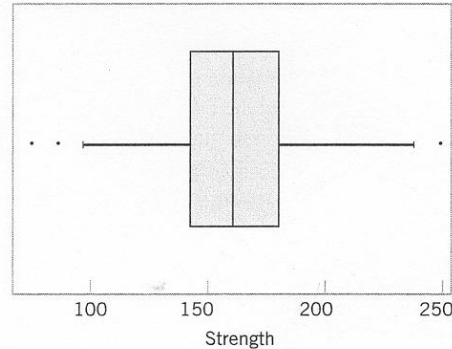


Figure 2-12 Box plot for compressive strength data in Table 2-2.

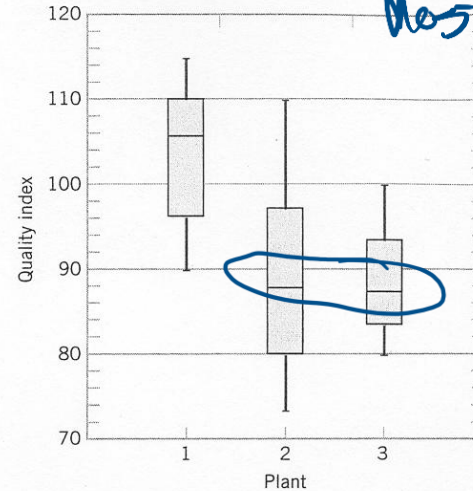


Figure 2-13 Comparative box plots of a quality index at three plants.

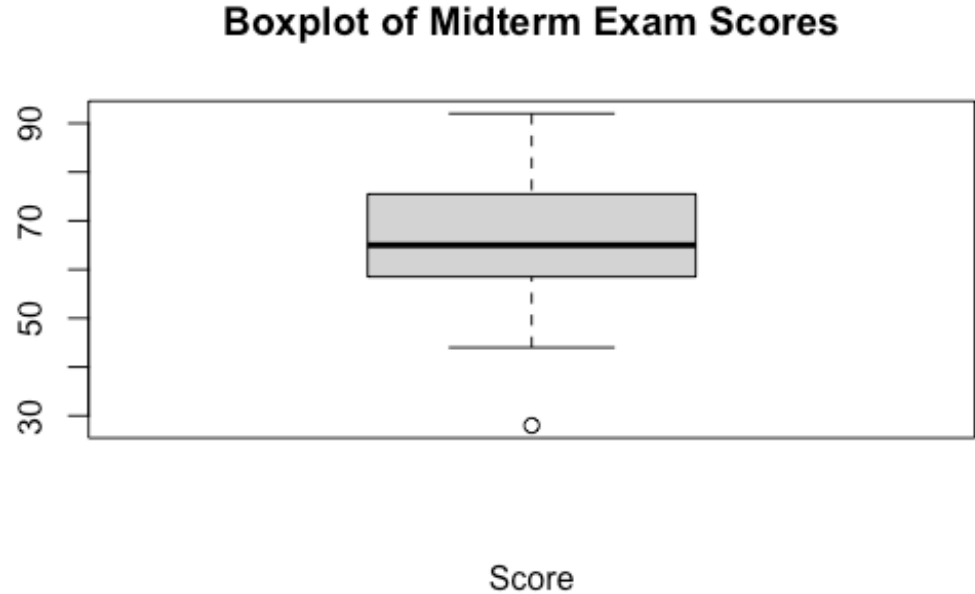
examples of boxplots

### Box Plot 3

R code for Box Plot

```
boxplot(midterm$MidtermExam,xlab='Score',main='Boxplot of Midterm Exam Scores')
```

R Box Plot output



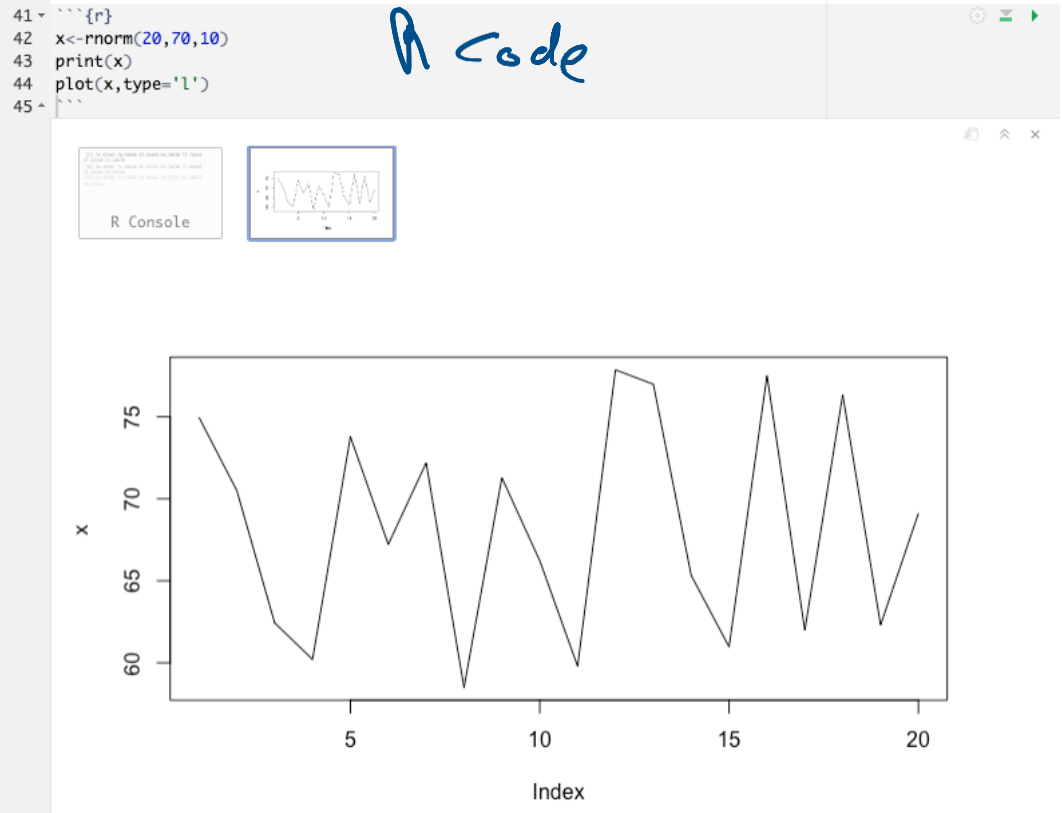
Boxplot of Midterm Exam Scores

# Time Series Plot

- A **time series plot** is a graph in which the vertical axis denotes the observed value of the variable (say  $x$ ) and the horizontal axis denotes time
- Excellent tool for detecting:
  - trends,
  - cycles,
  - other non-random patterns



## Time Series Plot in R



Time Series Plot

# Probability Plotting

- **Probability plotting** is a graphical method of determining whether sample data conform to a hypothesized distribution
- Used for validating assumptions
- Alternative to hypothesis testing

## Construction

1. Sort the data from smallest to largest, .
2.  $x_{(1)}, x_{(2)}, \dots, x_{(n)}$
3. Calculate the observed cumulative frequency  $(j - 0.5)/n$

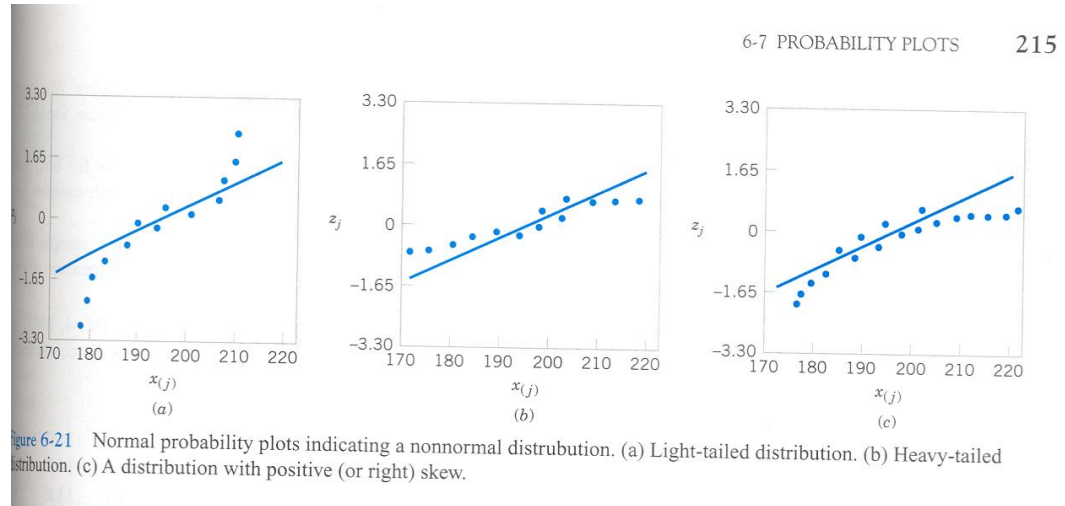
For the normal distribution find  $z_j$  that satisfies

$$\frac{j - 0.5}{n} = P(Z \leq z_j) = \Phi(z_j)$$

3. Plot  $z_j$  versus  $x_{(j)}$  on special graph paper

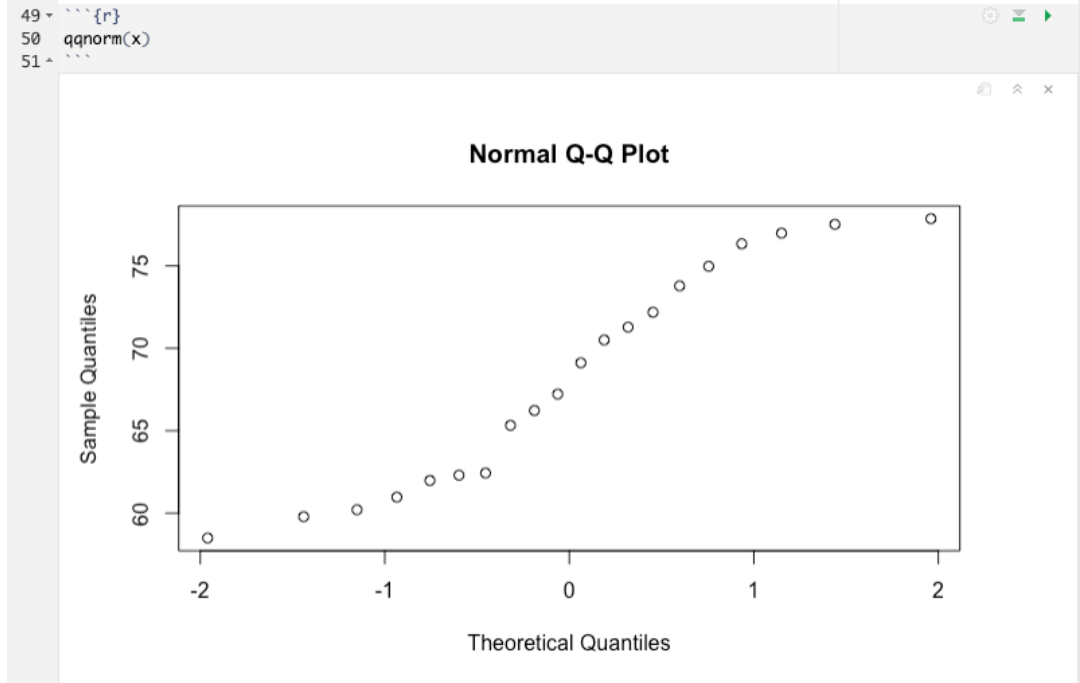
## Usage

If the data plots as a straight line, the assumed distribution is correct



normal probability plots from textbook, figure 6.21 on page  
215

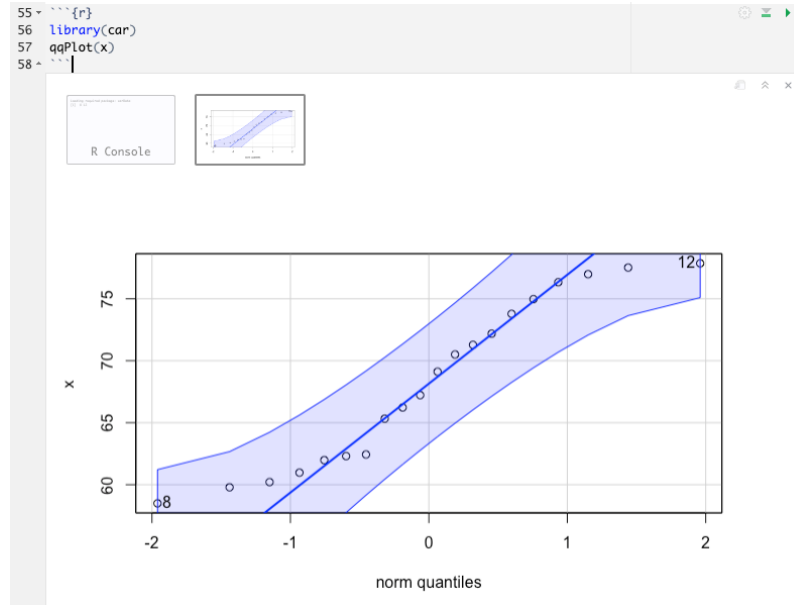
## Probability Plot Example 1 in R



Normal Probability Plot

# Probability Plot Example 2

- Difficulty from example one is how close to straight is “good enough”
- Add confidence bands to normal probability plot
  - Requires package car to be added to R
  - If all points are within the band, we are 95% confident that the sample is from a normal distribution. However if one or more points are not within band, the data is not from a normal distribution

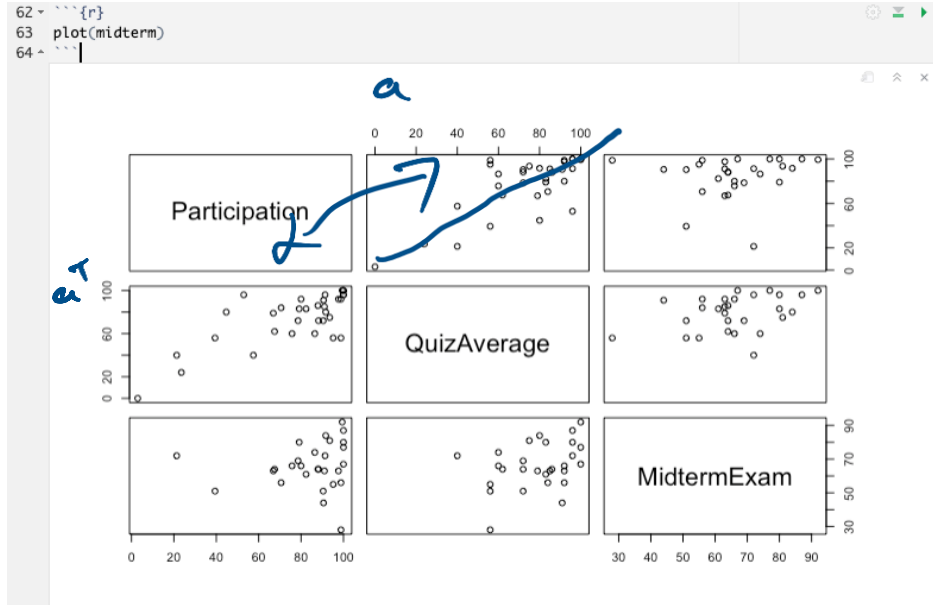


QQ Plot with band

# MULTIVARIATE DATA



# Matrix of Scatter Plot in R



Scatter Plots

# Covariance in R

```
67 ~~~{r}  
68 midterm_NA <- na.omit(midterm)  
69 print(cov(midterm_NA))  
70 ~~~
```

	Participation	QuizAverage	MidtermExam
Participation	340.16778	193.7847	28.75699
QuizAverage	193.78474	269.0899	81.17460
MidtermExam	28.75699	81.1746	188.43915

Covariance Matrix

# Correlation

```
74 ~~~{r}  
75 print(cor(midterm_NA))  
76 ~~~
```

	Participation	QuizAverage	MidtermExam
Participation	1.0000000	0.6405076	0.1135825
QuizAverage	0.6405076	1.0000000	0.3604839
MidtermExam	0.1135825	0.3604839	1.0000000

## Correlation Matrix

# Chapter 7 Overview

- Chapter 7 contains a detailed explanation of point estimates for parameters
- Much of this chapter is of a highly statistical nature and will not be covered in this course
- Key concepts we will discuss are:
  - Statistical inference
  - Statistic
  - Sampling distribution
  - Point estimator
  - Unbiased estimate
  - MVUE estimator
  - Central limit theorem
  - Sampling distributions

# Statistical Inference

- Montgomery gives the following description of statistical inference.

The field of statistical inference consists of those methods used to make decisions or to draw conclusions about a population. These methods utilize the information contained in a sample from the population in drawing conclusions. This chapter begins our study of the statistical methods used for inference and decision making.

- Statistical inference may be divided into two major areas:  
parameter estimation and hypothesis testing

Chapter 8

Chapter 9 & 11

## Point Estimate

- Montgomery states that “In practice, the engineer will use sample data to compute a number that is in some sense a reasonable value (or guess) of the true mean. This number is called a **point estimate**.”
- Discuss examples
- A formal definition of a point estimate is  
A **point estimate** of some population parameter  $\theta$  is a single numerical value  $\hat{\theta}$  of a statistic  $\hat{\Theta}$ . The statistic  $\hat{\Theta}$  is called the point estimate.
- Notice the use of the “hat” notation to denote a point estimate

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \approx f(x_1, x_2, \dots, x_n)$$

## Statistic

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{as } n \rightarrow \infty \quad \bar{X} \text{ has a normal distribution}$$

- Point estimate requires a sample of random observations, say  $X_1, X_2, \dots, X_n$
- Any function of the sampled random variables is called a statistic
- The function of the random variables is itself a random variable
- Thus, the sample mean  $\bar{x}$  and the sample variance  $s^2$  are both statistics and random variables

## **Properties of point estimators**

- We would like point estimates to be both accurate and precise
- An unbiased estimator addresses the accuracy criteria
- A minimum variance unbiased estimator addresses the precision criteria



## Unbiased Estimator

- The point estimator  $\hat{\Theta}$  is an **unbiased estimator** for the parameter  $\theta$  if

$$E(\hat{\Theta}) = \theta$$

- If the point estimator is not unbiased, then the difference

$$E(\hat{\Theta}) - \theta$$

is called the **bias** of the estimator  $\hat{\Theta}$

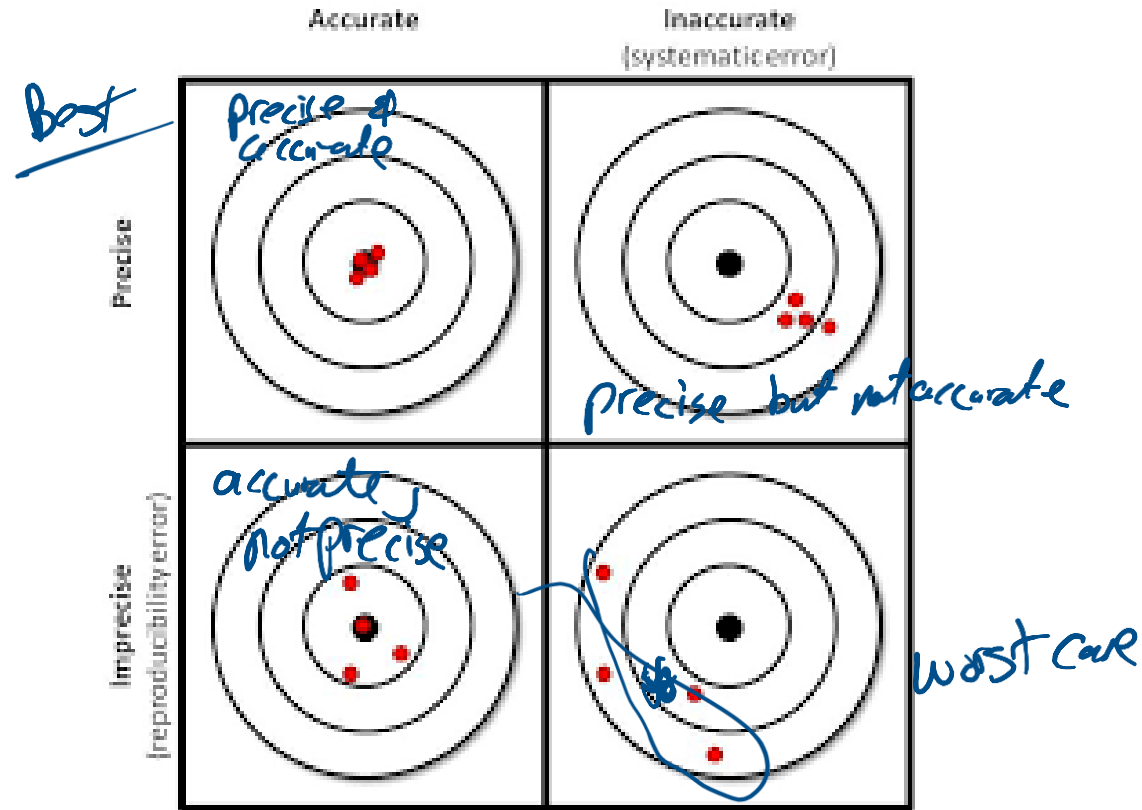
## MVUE

- Montgomery gives the following definition of a minimum variance unbiased estimator (MVUE)

If we consider all unbiased estimators of  $\theta$ , the one with the smallest variance is called the minimum variance unbiased estimator

- An important fact is that the sample mean  $\bar{x}$  is the MVUE for  $\mu$  when the data comes from a normal distribution

## Accuracy vs. Precision



graph of accuracy vs. precision

## **Sampling Distribution**

- The probability distribution of a statistic is called a **sampling distribution**

## Central Limit Theorem

- Definition of the Central Limit Theorem is

If  $X_1, X_2, \dots, X_n$  is a random sample of size  $n$  taken from a population (either finite or infinite) with mean  $\mu$  and finite variance  $\sigma^2$ , and if  $\bar{X}$  is the sample mean, the limiting form of the distribution of

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

as  $n \rightarrow \infty$ , is the standard normal distribution

- Important result because for sufficiently large  $n$ , the sampling distribution of  $\bar{X}$  is normally distribution
- This is a fundamental result that will be used extensively in the next four chapters of the textbook.